

Response to Professor Stephen Jenkins' comments on the World Income Inequality Database (WIID)

Nina Badgaiyan¹ · Jukka Pirttilä¹ · Finn Tarp^{1,2}

Received: 2 June 2015 / Accepted: 3 June 2015 / Published online: 16 June 2015
© Springer Science+Business Media New York 2015

Professor Stephen Jenkins (*this issue*) has conducted an extremely careful and insightful analysis of two datasets, WIID and SWIID. In this short response, we focus on his review of the WIID, maintained and published by UNU-WIDER in agreement with the World Bank. We wish to highlight at the outset that we are very grateful for Professor Jenkins' expert advice and suggestions; and note that the WIID has over the past couple of years been developed in the direction recommended in the appraisal.

This response provides, first, some background to the development and basic philosophy of the WIID, with particular attention to data comparability and proper documentation. We subsequently outline the recent changes made to the database, now available on the internet as the revised version WIID3.0b.¹ Finally, we offer some brief concluding remarks.

1 The WIID – background and principles

The WIID was initially assembled in 1997–99 for the UNU-WIDER-UNDP project 'Rising Income Inequality and Poverty Reduction: Are They Compatible?' directed by Professor Giovanni Andrea Cornia, the then Director of UNU-WIDER. The original WIID incorporated the Deininger and Squire (1996) data set, increased the country coverage, extended the timeframe, augmented the number of distributional indicators, expanded the number of data observations and a process of carefully assessing the quality of the data was initiated. As more observations were added UNU-WIDER decided in consultation with the World Bank to make the database publicly available. The objective was to facilitate further analysis and debate on global inequality. This resulted in WIID1.0, published in September 2000.

¹ See http://www.wider.unu.edu/research/WIID3-0B/en_GB/database/

✉ Finn Tarp
Finn.Tarp@econ.ku.dk

¹ UNU-WIDER, Helsinki, Finland

² University of Copenhagen, Copenhagen, Denmark

The WIID database was subsequently updated and improved as part of the UNU-WIDER project 'Global Trends in Inequality and Poverty' led by Professor Tony Shorrocks. He was the UNU-WIDER Director from 2001-09; and the second, substantially improved version of the WIID was released in June 2005. In the following years WIDER continued to update the WIID, allowing further improvements in the quality of data, the extent of international coverage, and the length of the time series. This process produced WIID2.0c, dated 2008, which is the version Professor Jenkins appraised.

The philosophy underlying the preparation of WIID2 was to gather income inequality indices (the Gini index and quintile shares) for as many countries and years as possible. This was done being cognizant of the fact that observations in a secondary database such as the WIID will originate from different sources and refer to a variety of income and population concepts, sample sizes and statistical methods. To be as transparent and helpful to the users of WIID as possible, the UNU-WIDER team did its best to specify the conceptual base for each observation and to otherwise document the data well. One of the requirements of Atkinson and Brandolini (2001) (referred to by Professor Jenkins) was, indeed, to provide a thorough documentation of the data.

Accordingly, WIID regularly provides a large number of Gini observations for each country-year pair. Observations may refer to different concepts of inequality (for example, whether the Gini is for market income, disposable income or consumption inequality). WIID values may have been calculated on the basis of the full population or using a subset (rural vs urban) only, or they may rely on different equivalence scales. This means that there is no single 'correct' WIID value for a country in a particular year. This implies that when cross-country comparisons of income inequality are made by analysts, the series relied on must be carefully selected so they are comparable. Sometimes this is not possible: for example for the majority of developing countries, all Gini values are calculated on the basis of consumption, and these should not be compared with market/gross income inequality in developed countries. We stress that this is not a shortcoming of WIID per se. It reflects the underlying data availability.

Importantly, the WIID also contains a quality classification of observations. The highest quality observations are not always available for all countries and years. Consequently, users of the database confront a trade-off between coverage and quality, which imply that conscious choices must be made.

It was by design that WIID contains all the necessary qualifying variables and a full documentation of the sources of data. The idea has throughout been to enable users to make sensible selections for the data they use and minimize the risk of comparing 'apples with oranges' in cross-country analyses. On this background, we are encouraged that WIID has – in Professor Jenkins' assessment – 'successfully implemented Atkinson and Brandolini's (2001) recommendations regarding *construction and development* of secondary databases'.

One consequence of the differences in methodology in data-gathering when constructing inequality indices in different countries is that many countries do not have data for all the years and all the notions of inequality the analyst would wish to have. The strategy in developing the WIID is to provide the user with data actually available, and avoid filling in major gaps by imputations or any other means. Imputations may be reasonable when filling in a missing observation in a particular time-series for an individual country where data is available for other years. They do, however, become problematic if, e.g. Gini values for net income are imputed on the basis of information from other countries, when limited or no information is available for the country in question for that income concept.

We find that gathering and presenting the original data in a clear and transparent manner is the only credible way forward, and we agree with Professor Jenkins that researchers need

to come fully to grips with the actual, underlying data with which they work. If imputations are added, they should be provided via a separate, auxiliary dataset, so the changes made to the original data can be highlighted and understood.

We also agree with the recommendation by Professor Jenkins that users need to be explicit and transparent in the choices they make when they use WIID. The documentation of WIID2 explicitly encouraged the users to do so; and we have further emphasized this point in the documentation of the latest version of the data, WIID3.0b.

2 WIID3: changes made to the database following Professor Jenkins' comments

As already noted, we have taken on-board Professor Jenkins' comments and suggestions on WIID, version 2.0c when compiling WIID3.0b (dated September 2014). In addition, WIID3.0b includes a range of new observations including data for 2012.

It may be helpful to provide a brief summary of the changes made to WIID2. The changes are related, first, to how to make the database more user-friendly, and, second, how to deal with the treatment of multiple observations.

All changes suggested by Professor Jenkins for making the database more accessible to the user have been incorporated. Thus (i) country documentation for the UK and Vietnam have been added; (ii) data is now sorted by country year (1st para); (iii) the currency reference variable has been divided into two variables: currency unit and reference period; (iv) misspellings and typographical errors have been corrected, and redundant decimal points from Reported Gini variable have been removed (1st para); (v) all numeric variables (1st para) including the quality variable are now value-labelled; (vi) the year variable is now numeric with mid points being used for observations such as 1953–55 (1st para), and the redundant variable AK has been dropped (1st para); (vii) two-letter country identifiers have been provided (2nd para) and country identifiers have been assigned unique numbers (2nd para); (viii) country documentation for all countries is now in one single pdf file (2nd para); (ix) Gini variable has been dropped since it is highly correlated with Reported Gini variable (2nd para); and (x) as far as possible, the country documentation and variable documentation have been closely integrated (2nd para).

Regarding multiple observations, Professor Jenkins suggests that they should all be dropped. We have done so in a number of cases. Thus, for instance, where there was overlapping information from EUROSTAT, SEDLAC and LIS only the latest information has been included. However, in some other cases we have retained multiple observations. This has been done in cases where we found that the multiple observations for the same country-year are in fact different in one or another dimension like for example source or area coverage. Also, we have not followed the recommendation on the Finnish series (Section 3.3, para 1) since a different equivalence scale (OECD, not OECD modified) is used. Similarly, we have not dropped observations merely because they are out of line with other observations for the same country-year though it was suggested that we do so, for instance in case of Hungary for 1990s etc. (1st para) or USA (Section 3.4). Instead, we use the quality variable referred to above to reflect our assessment of the relative usefulness of the observations. Likewise, we report both EUROSTAT and other WIID information for Ginis for EU countries giving the details of the two methodologies in the user guide section. The only exception to the rule regarding not dropping values merely because they are out of line, relates to the Jäntti series for Finland (2nd para), which we have dropped after consulting with the author.

3 Concluding remarks

We conclude by re-iterating that we fully support the basic message of Professor Jenkins: in empirical cross-country analysis of inequality, there is no easy way out from dealing with the details in making sure that the data selected are the best possible for the question in hand. Researchers need to make choices that are appropriate for the research question in focus. WIID aims to provide the information and documentation necessary to make such choices. Moreover, we share the view that researchers must communicate clearly and transparently the choices they have made for readers to fully understand their research strategy and for colleagues to be able to replicate results with minimum effort.

UNU-WIDER is committed to maintaining and updating the WIID and to facilitating research on inequality. In our further work, we plan to provide a user-friendly interface that will enable queries using the database and visualization. For this to work in a manageable way, the number of classes in some categorical variables (such as area and population coverage) will be aggregated to, say, three-five (in the case of area coverage, the natural classification is the entire country, urban vs rural and other). Even then, users must make the essential choices themselves, choosing the relevant series (be it market income or disposable income inequality) and the quality variance to be tolerated. At UNU-WIDER, we look forward to continuously engaging in policy relevant research on inequality in the years to come. We thank Professor Jenkins for an expert review and the editors for taking the initiative to put together the present issue of the *Journal of Economic Inequality*.

References

- Atkinson, A.B., Brandolini, A.: Promise and pitfalls in the use of secondary data-sets: Income inequality in OECD countries as a case study. *Journal of Economic Literature* **39**(3), 771–799 (2001)
- Deiningner, K., Squire, L.: A new data set measuring income inequality. *World Bank Economic Review* **10**(3), 565–591 (1996)